



Minería de datos y Patrones

Version 2024-I

Árboles de Decisión

[Capítulo 4]

Dr. José Ramón Iglesias

DSP-ASIC BUILDER GROUP

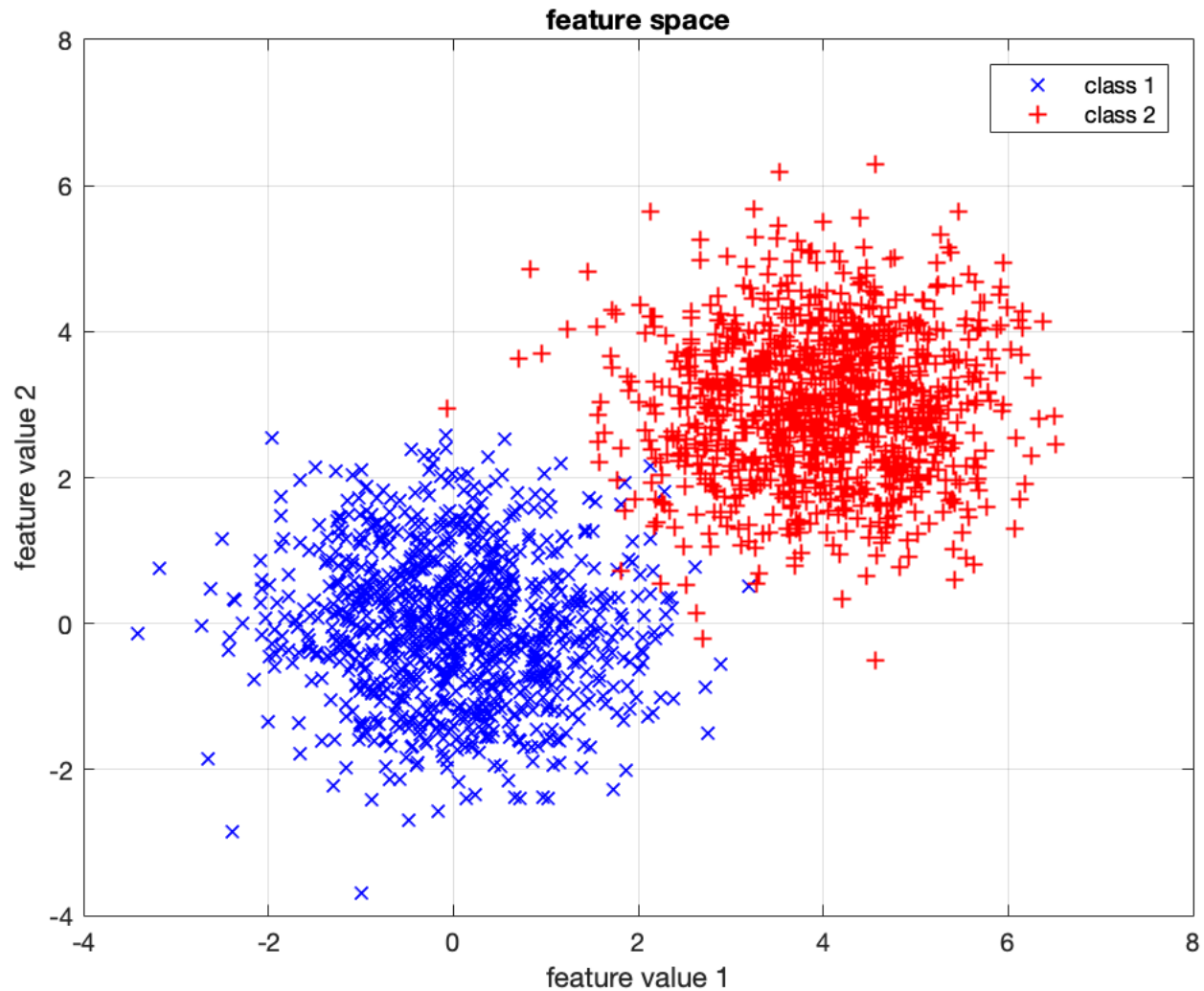
Director Semillero TRIAC

Ingeniería Electronica

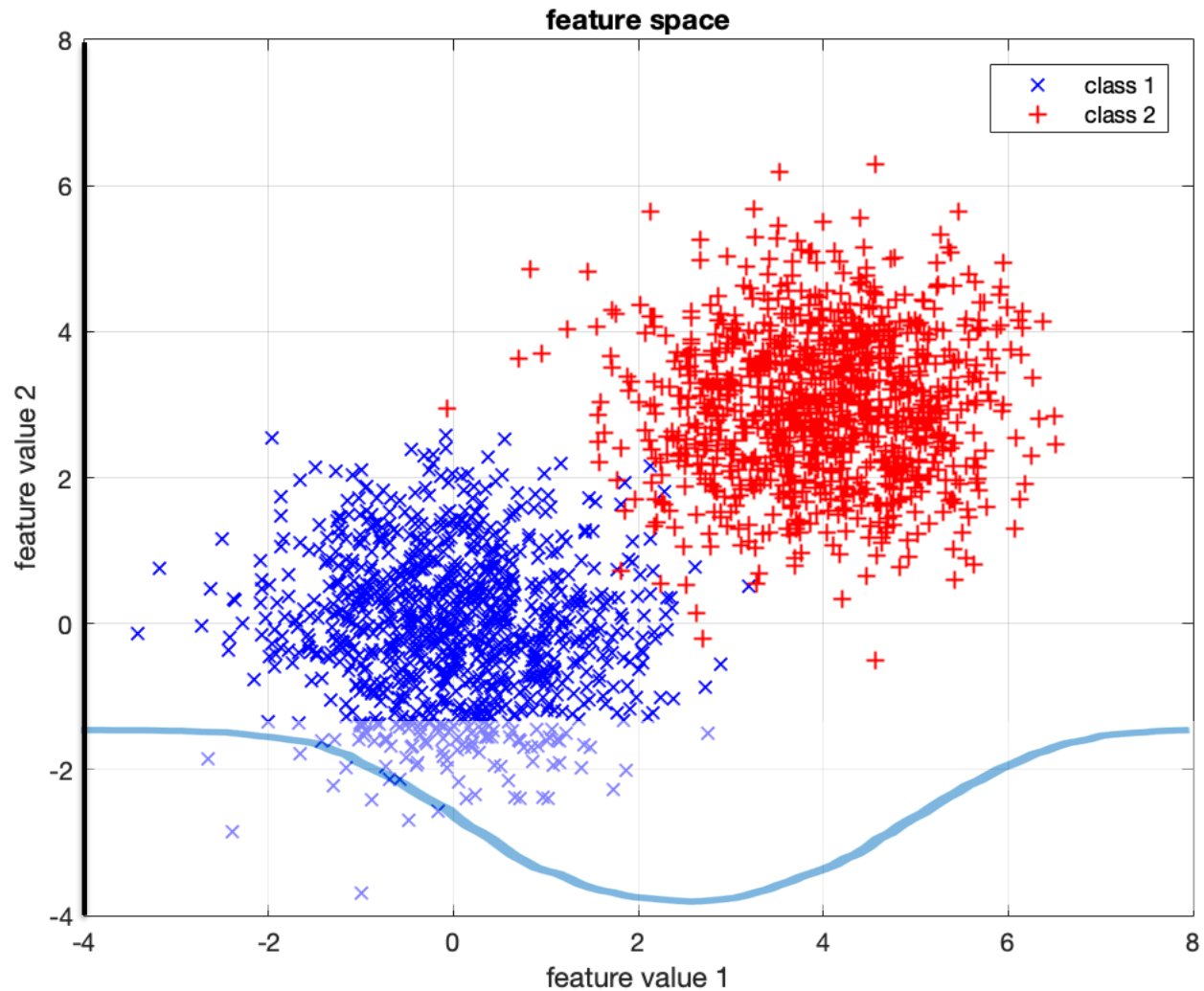
Universidad Popular del Cesar

Árboles de Decisión

Árboles de Decisión (Training data)

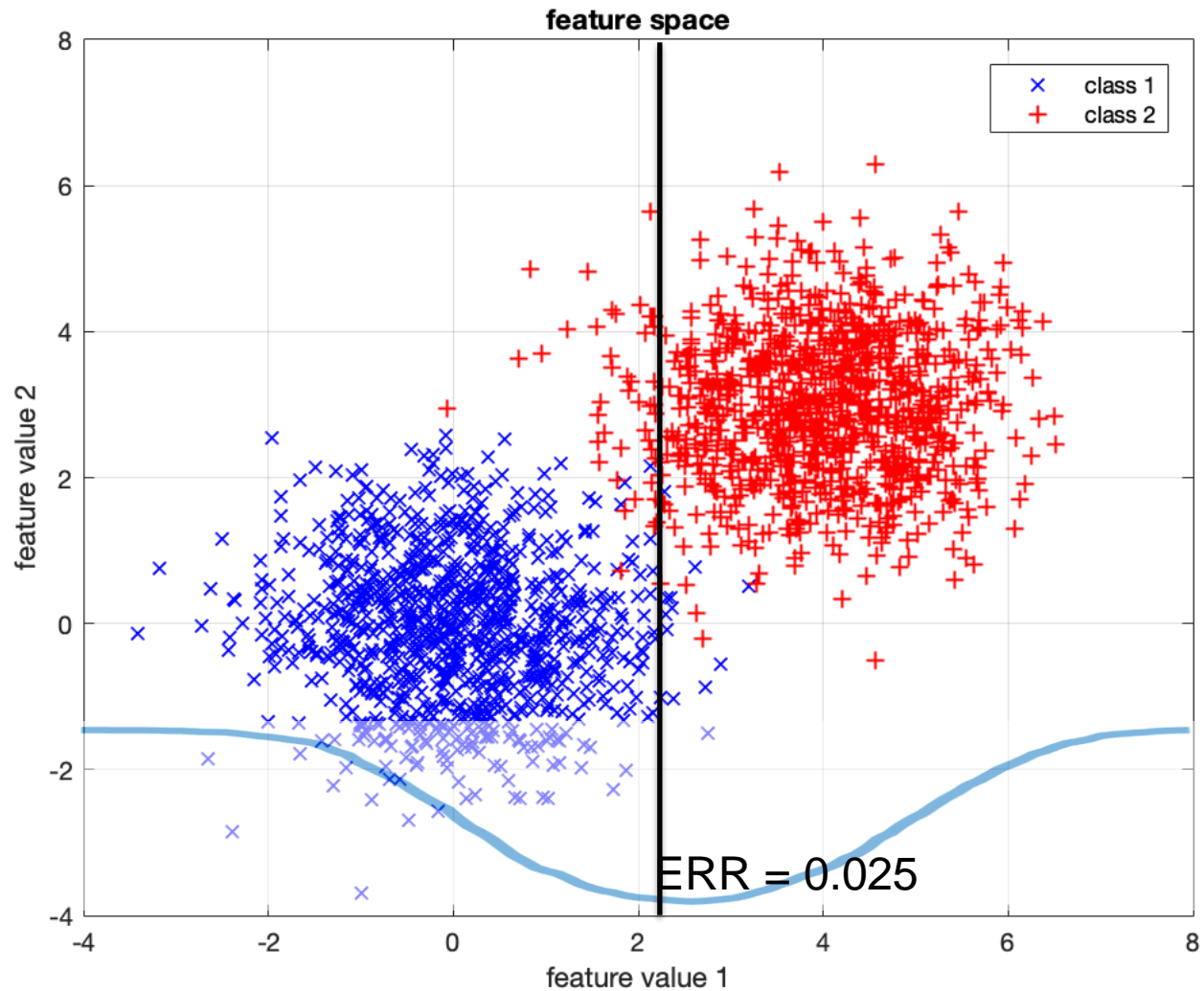


Árboles de Decisión (Training data)



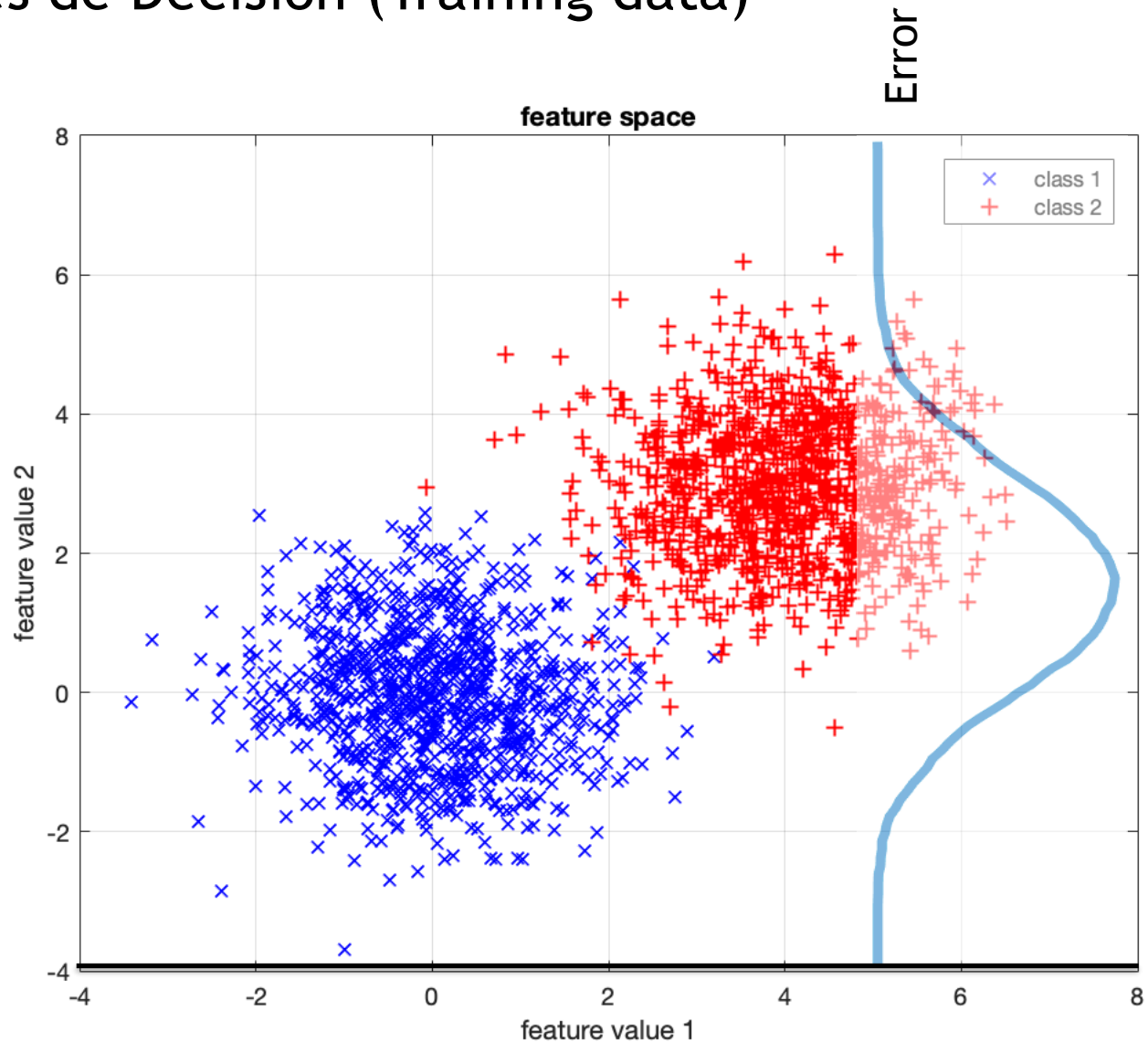
Error

Árboles de Decisión (Training data)

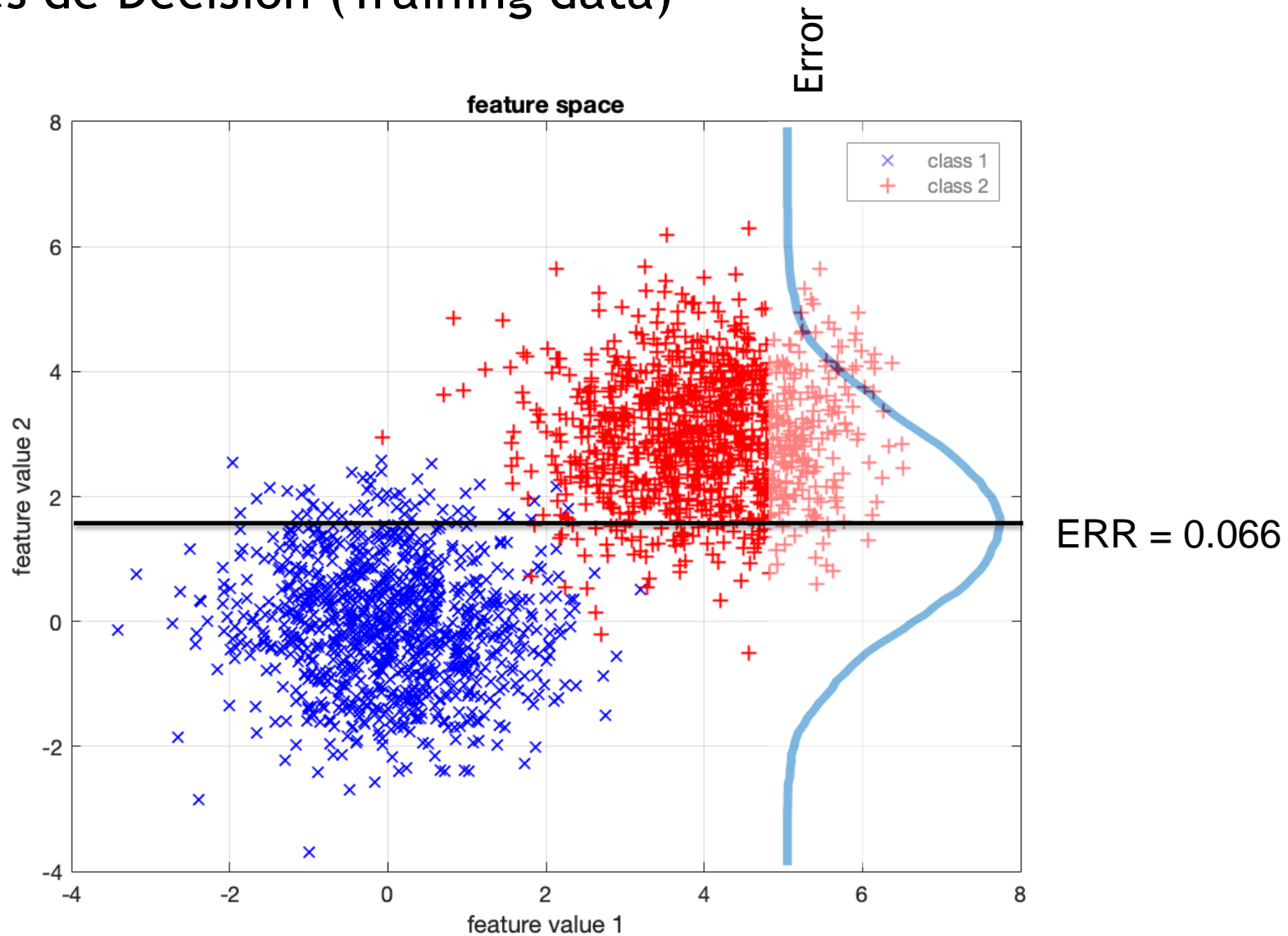


Error

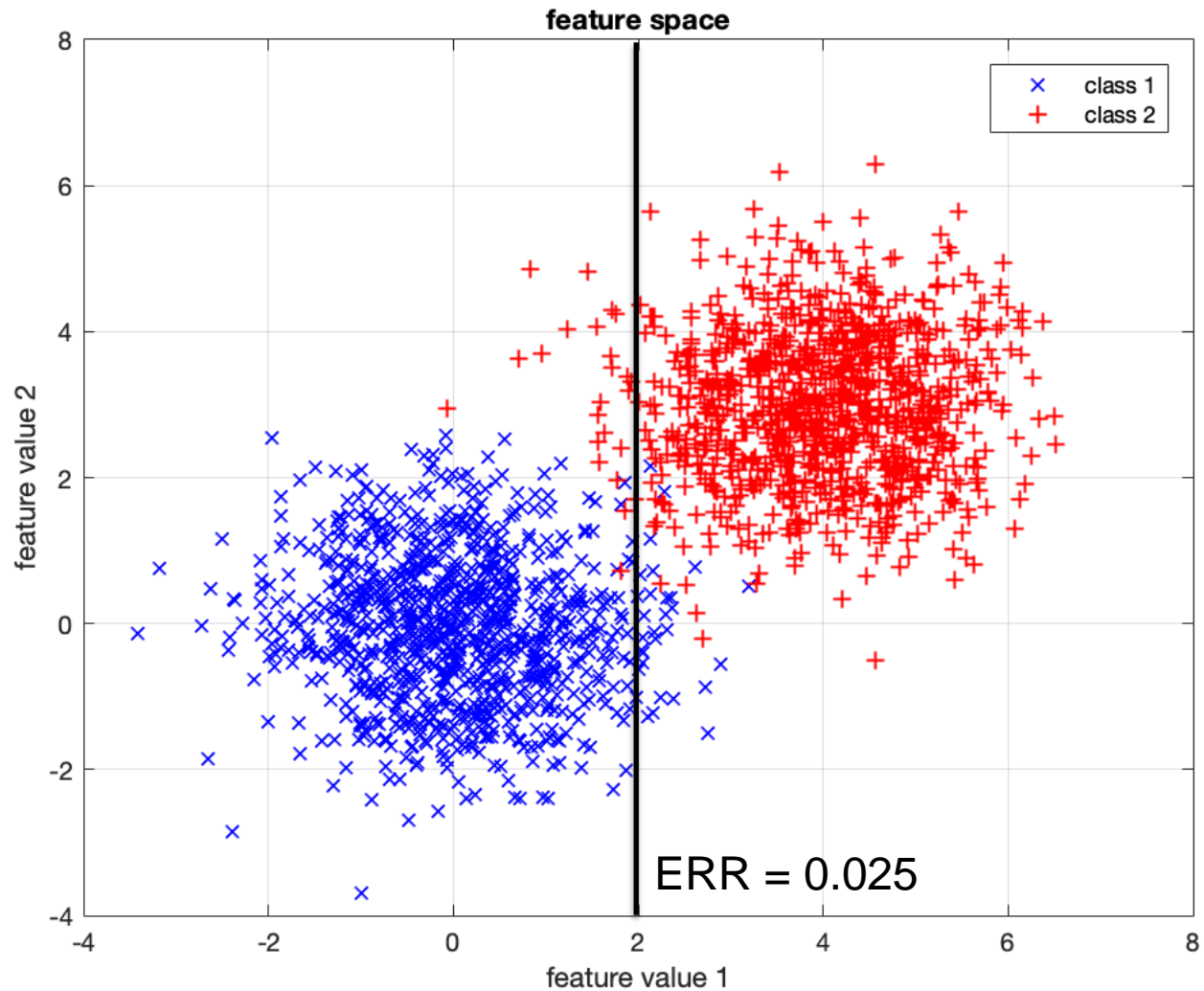
Árboles de Decisión (Training data)



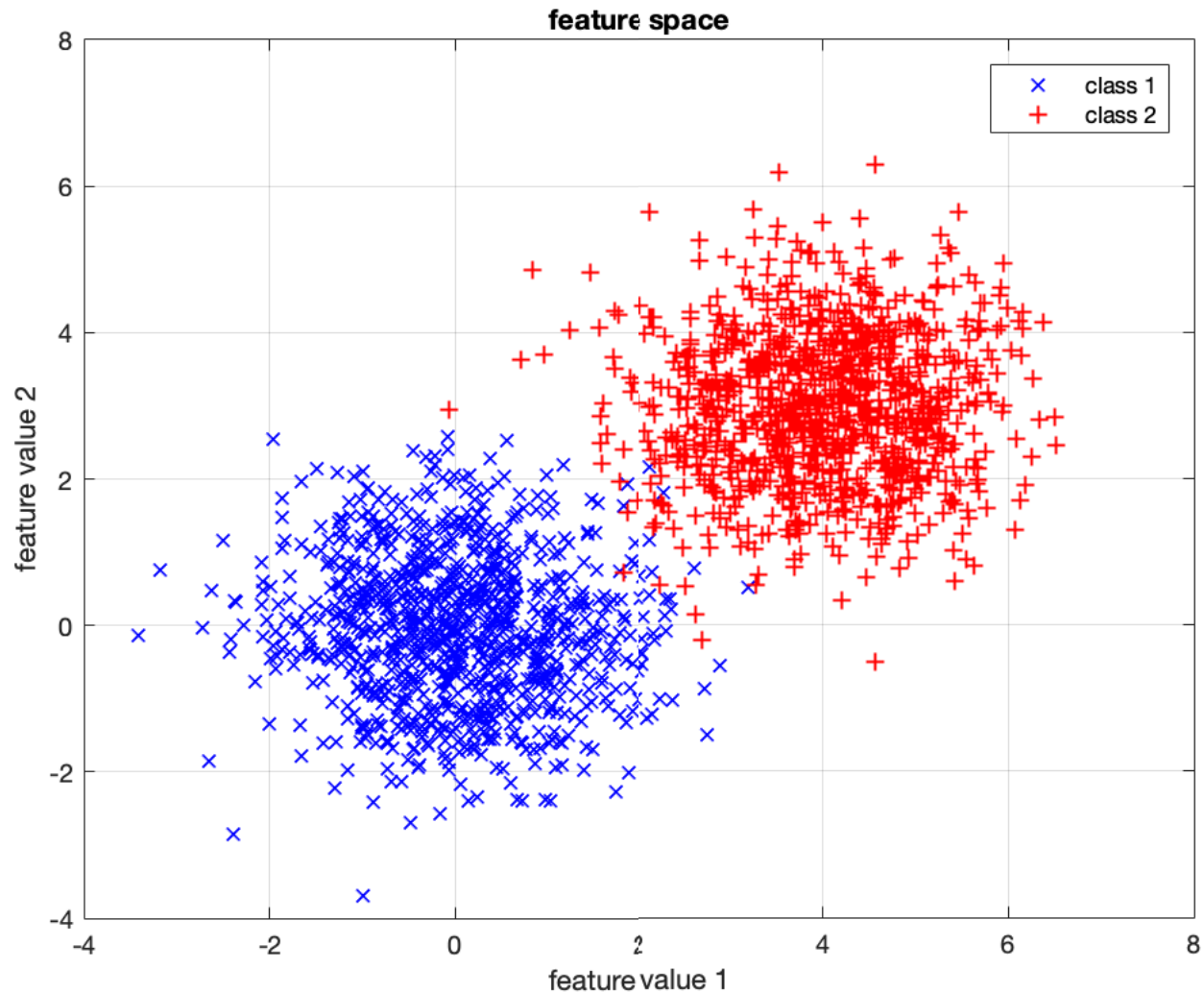
Árboles de Decisión (Training data)



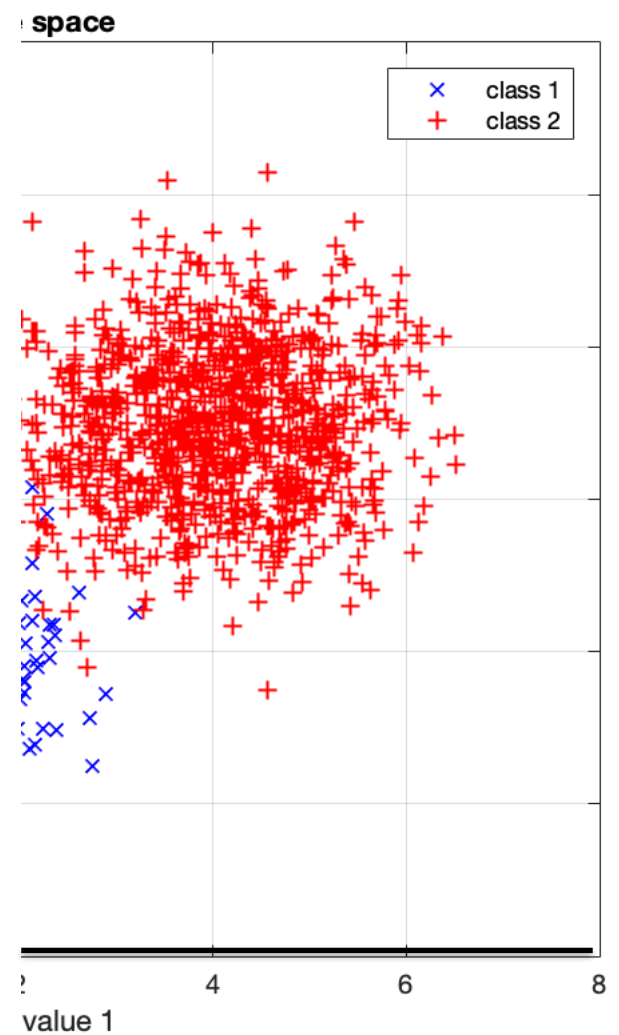
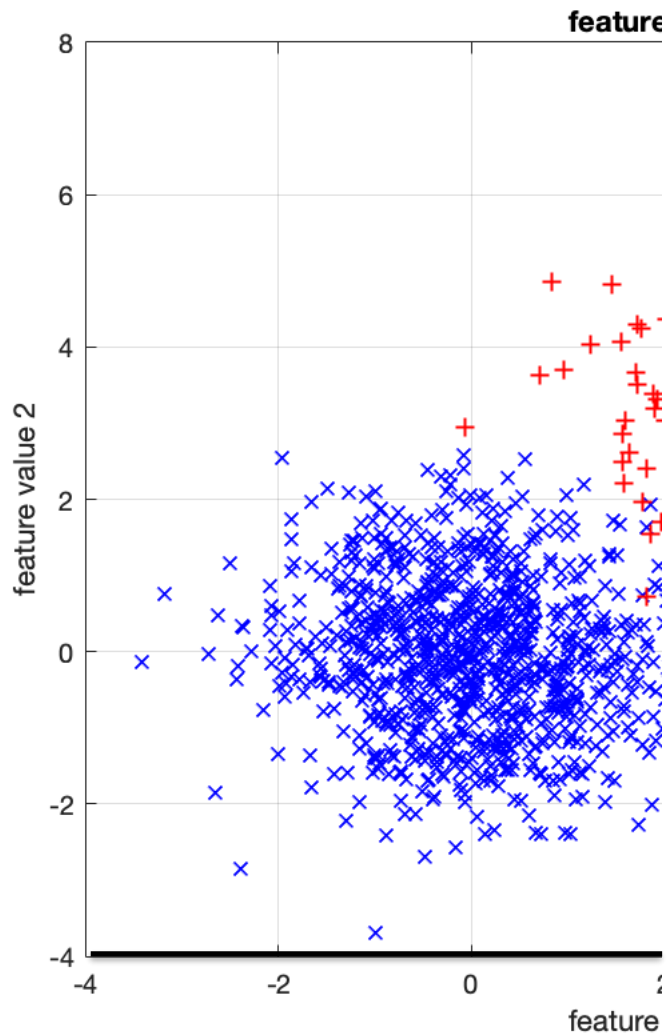
Árboles de Decisión (Training data)



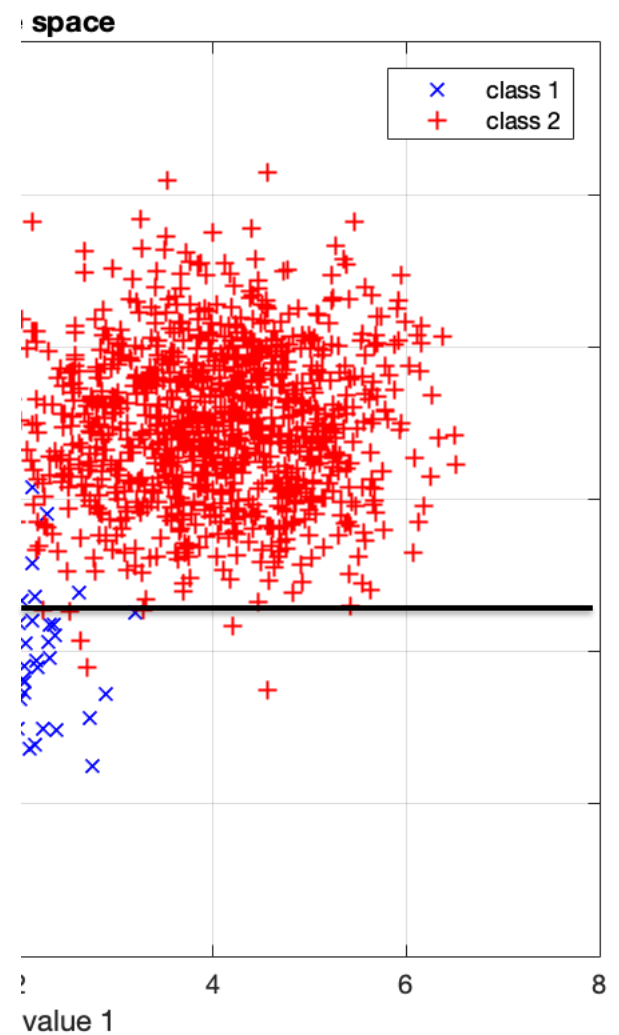
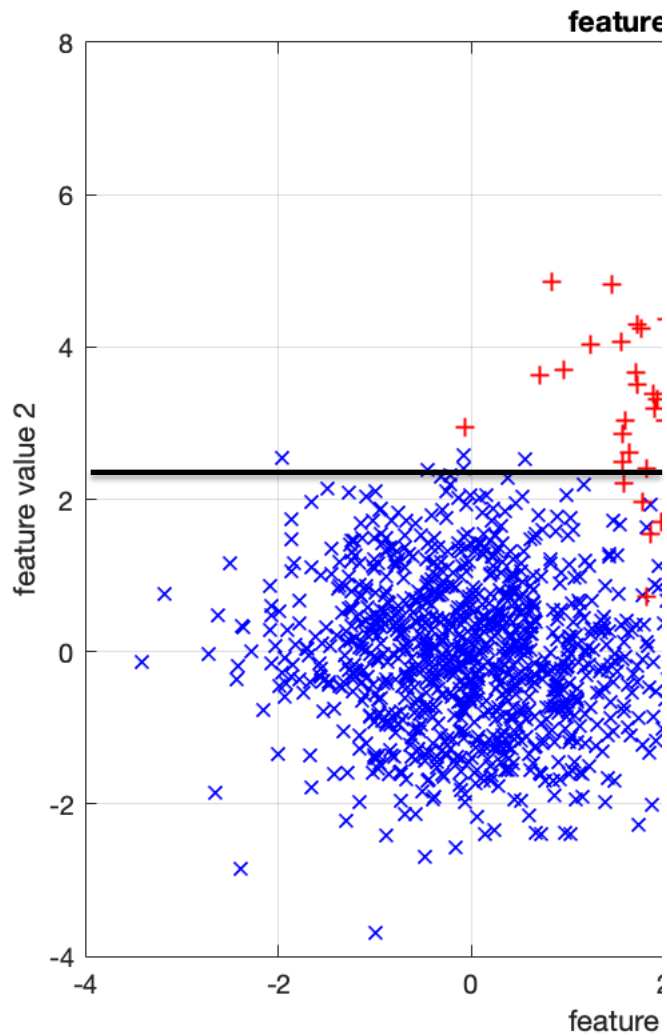
Árboles de Decisión (Training data)



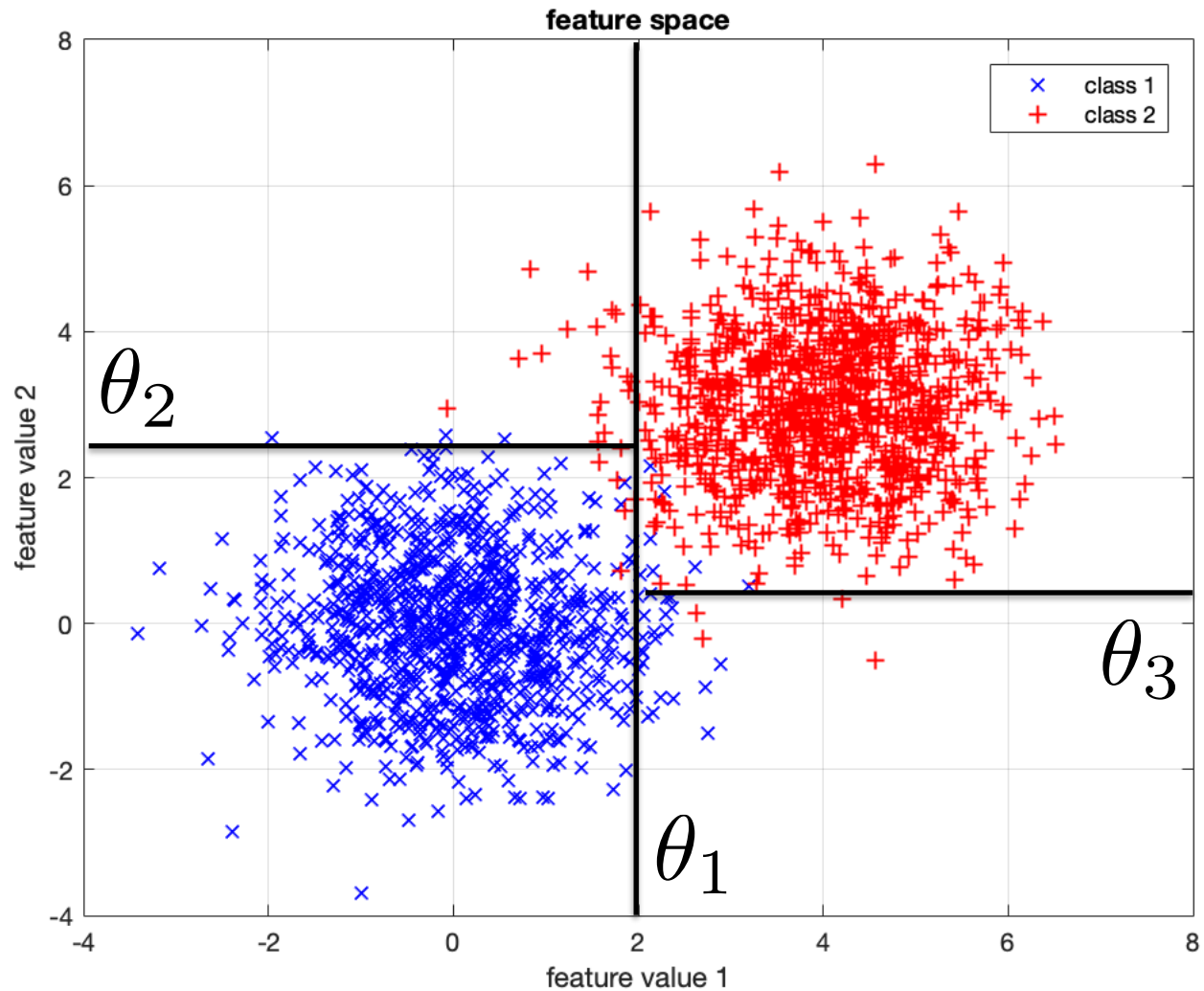
Árboles de Decisión (Training data)



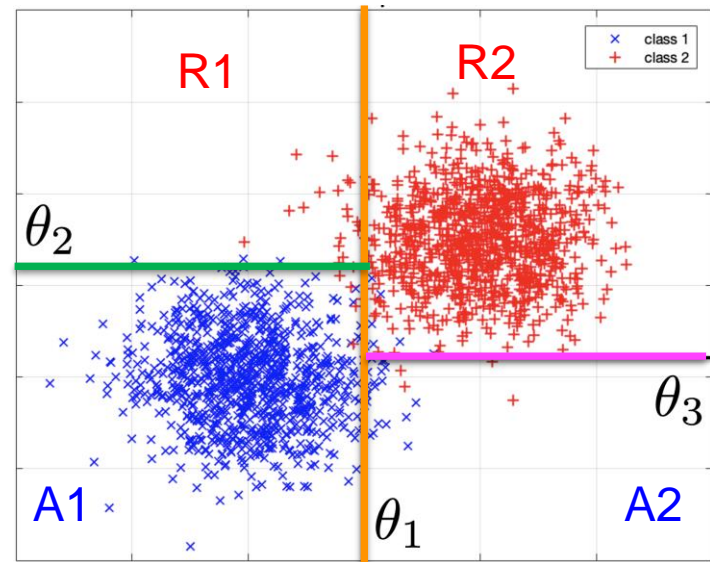
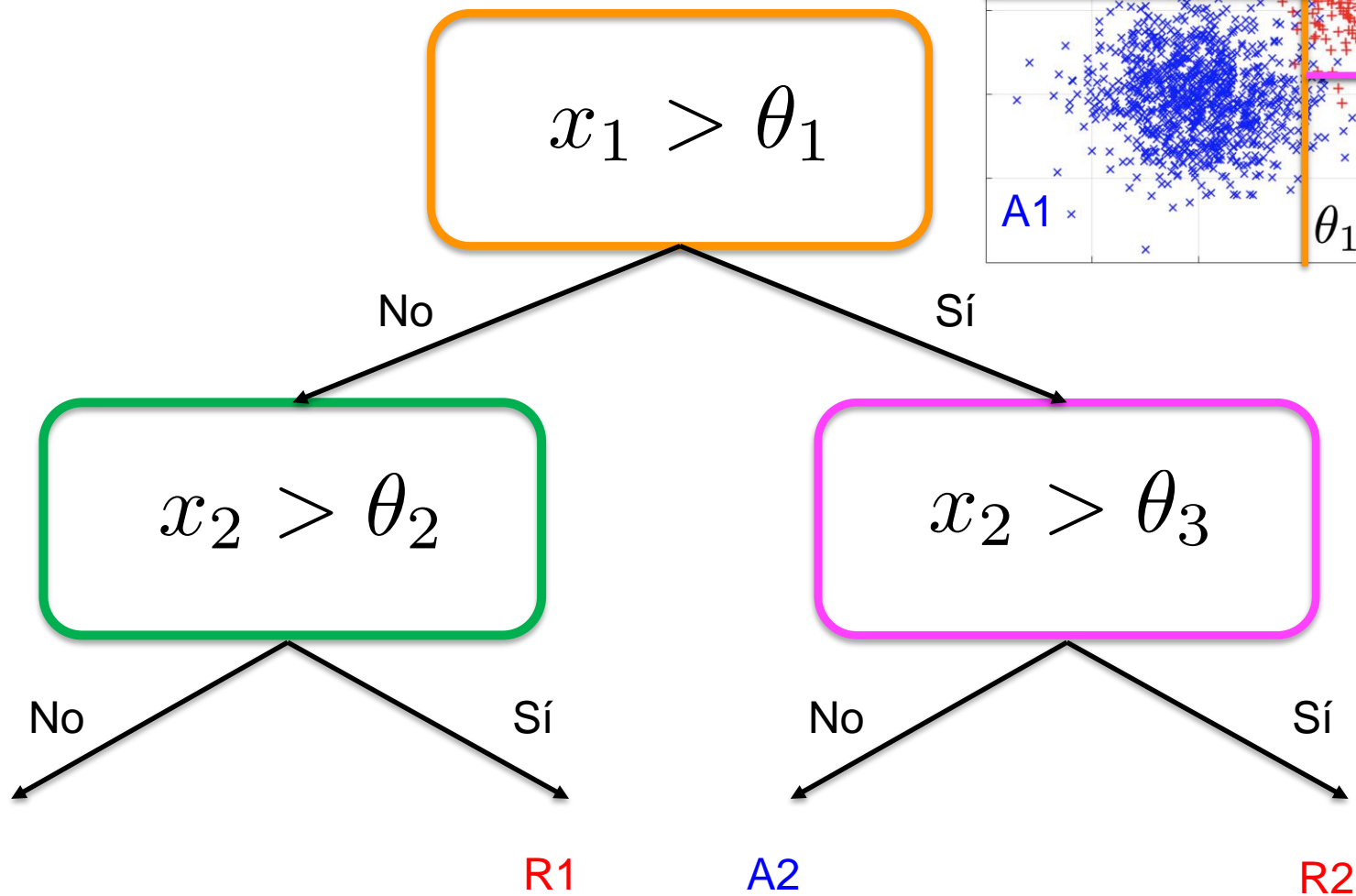
Árboles de Decisión (Training data)



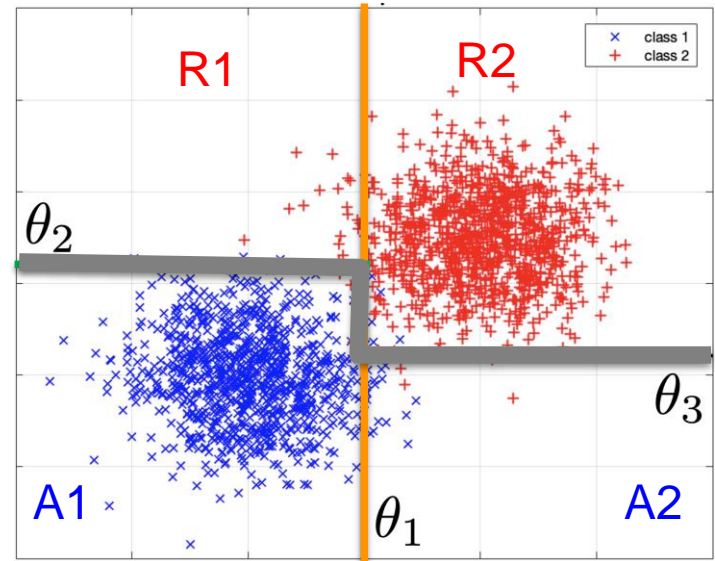
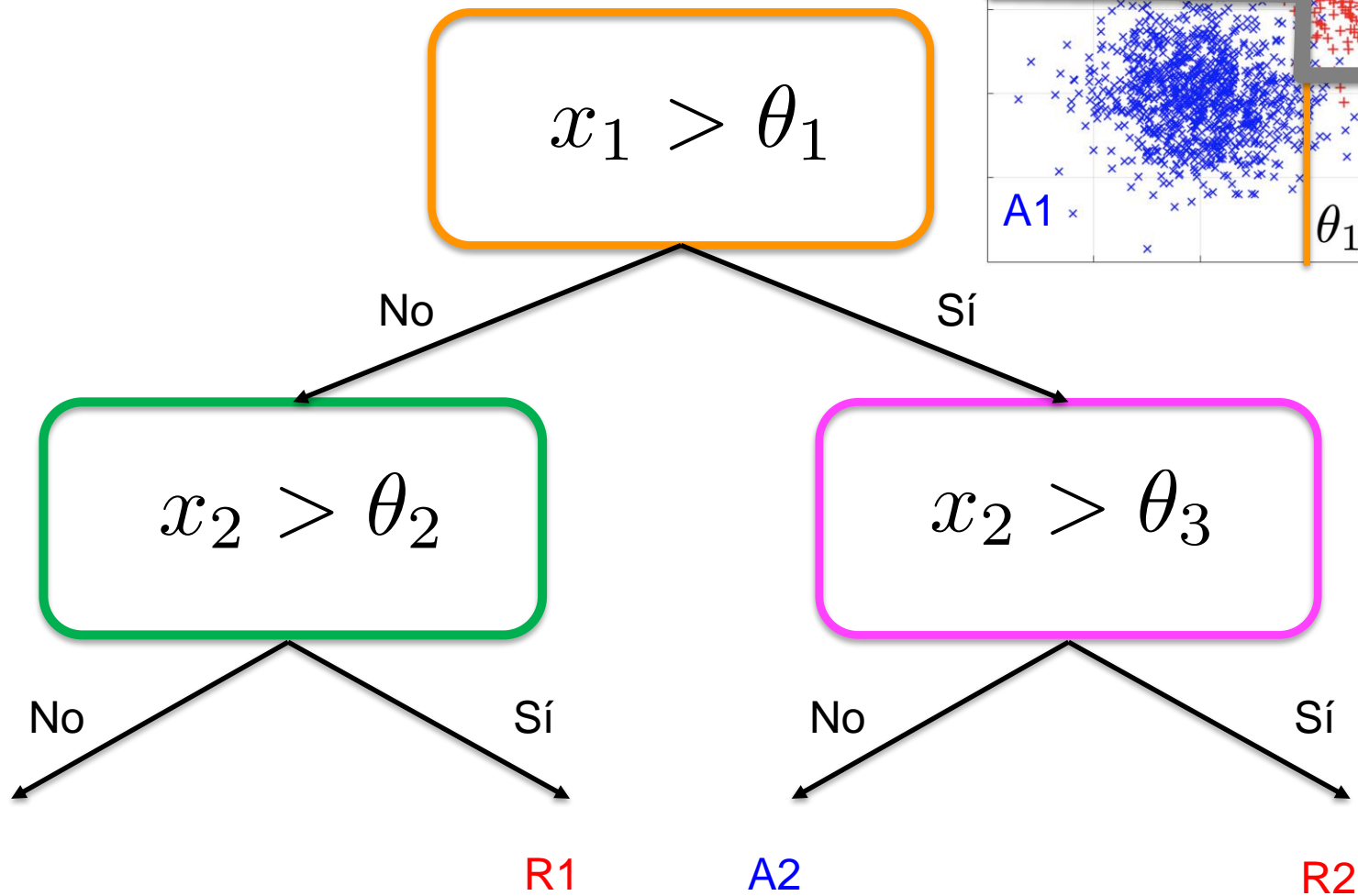
Árboles de Decisión (Training data)



Árbol



Árbol



Métricas usadas para el error:

- Error de clasificación 1 - Accuracy

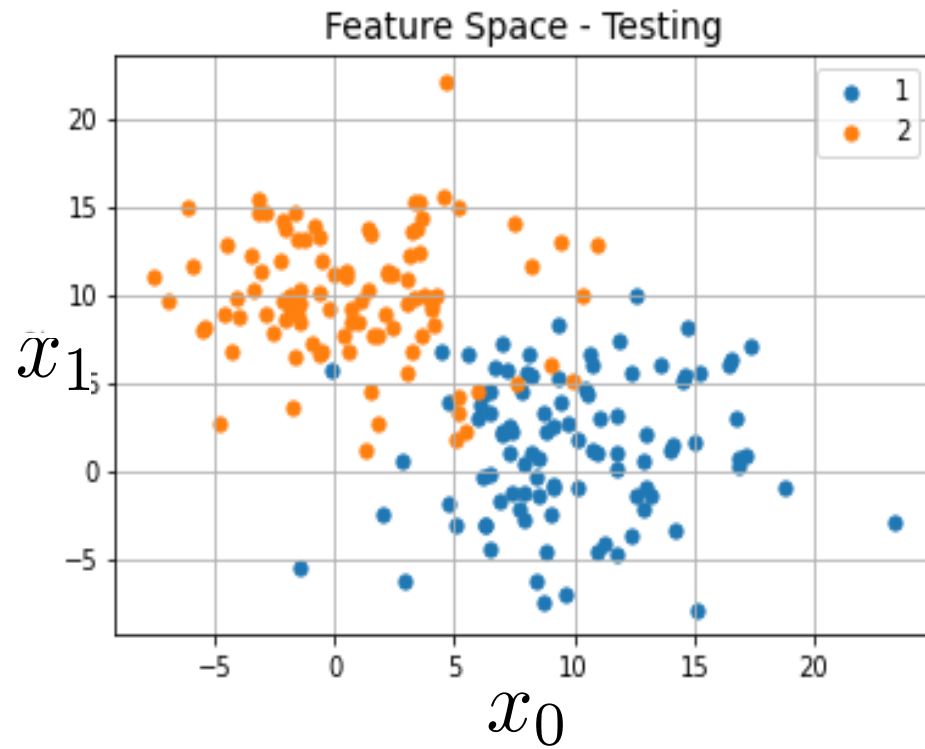
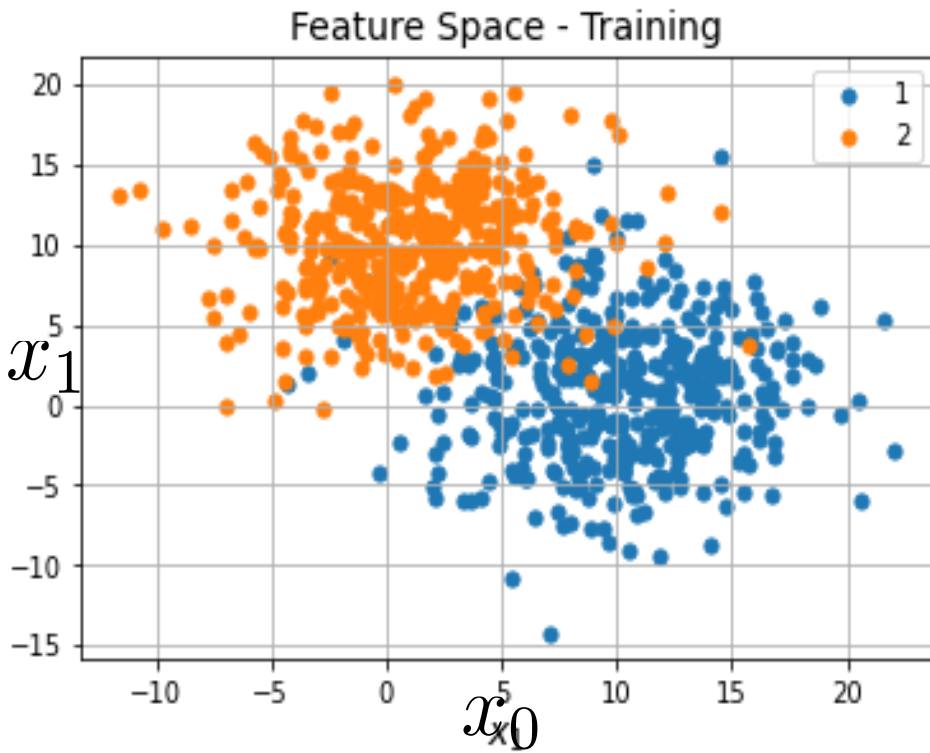
- Entropía $-\sum_{k=1}^K p_k \log_2(p_k)$

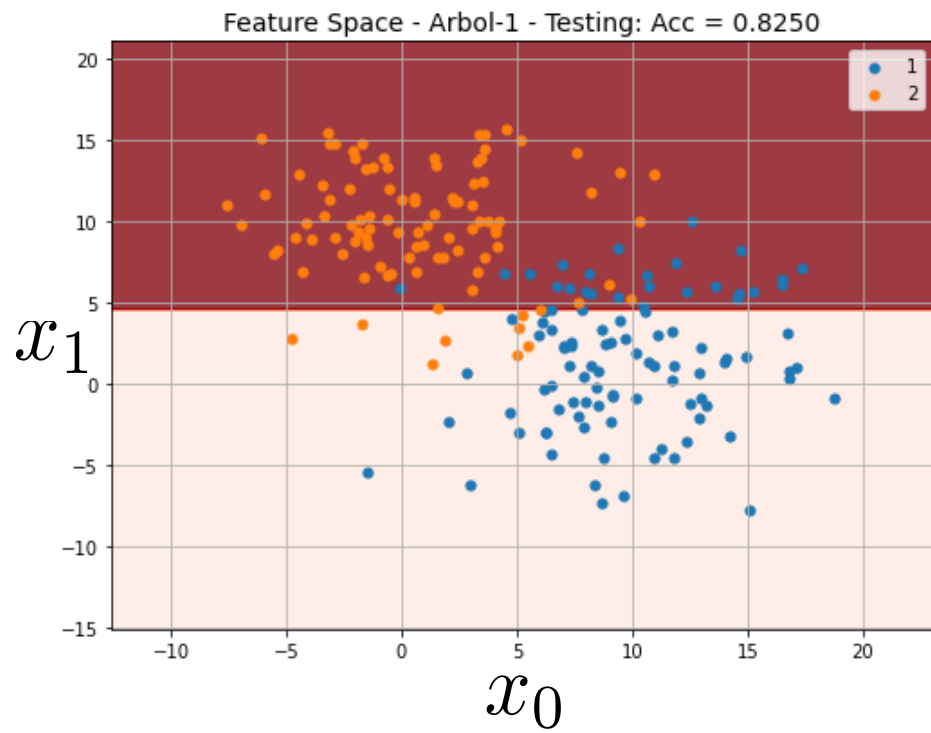
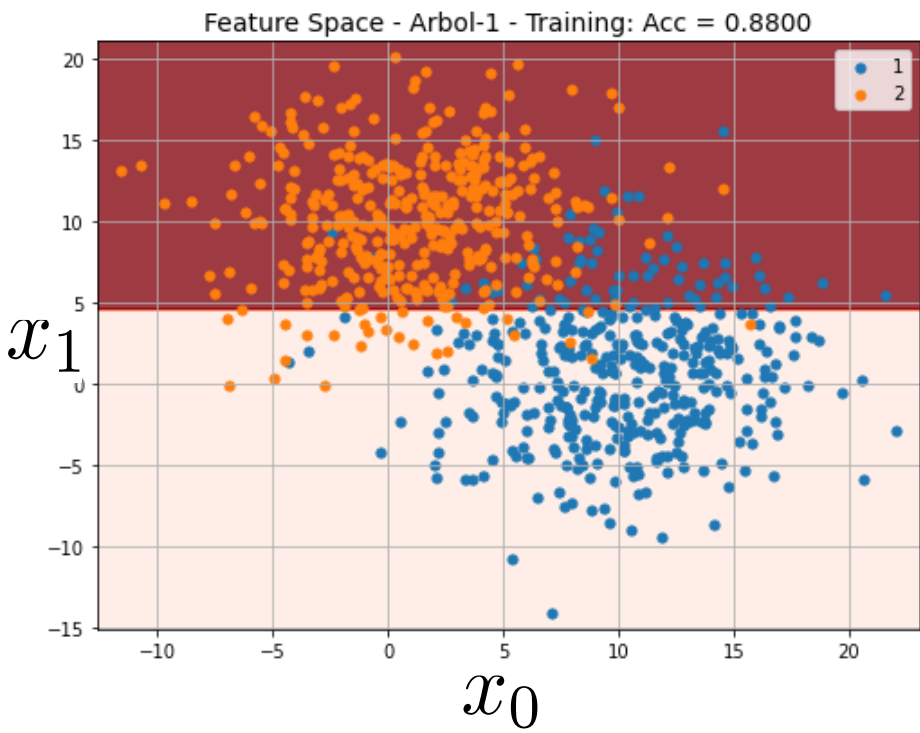
- Índice Gini $\sum_{k=1}^K p_k (1 - p_k)$

p_k Probabilidad de clasificar bien la clase k

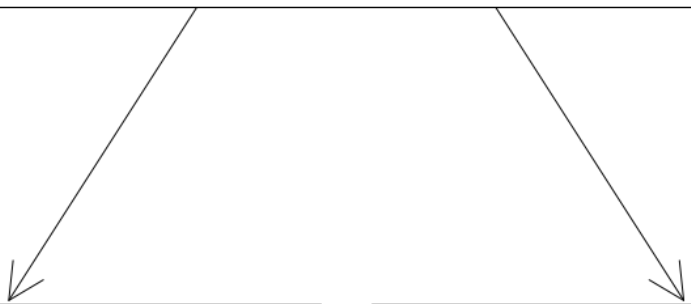
EJEMPLO

Datos Training/Testing



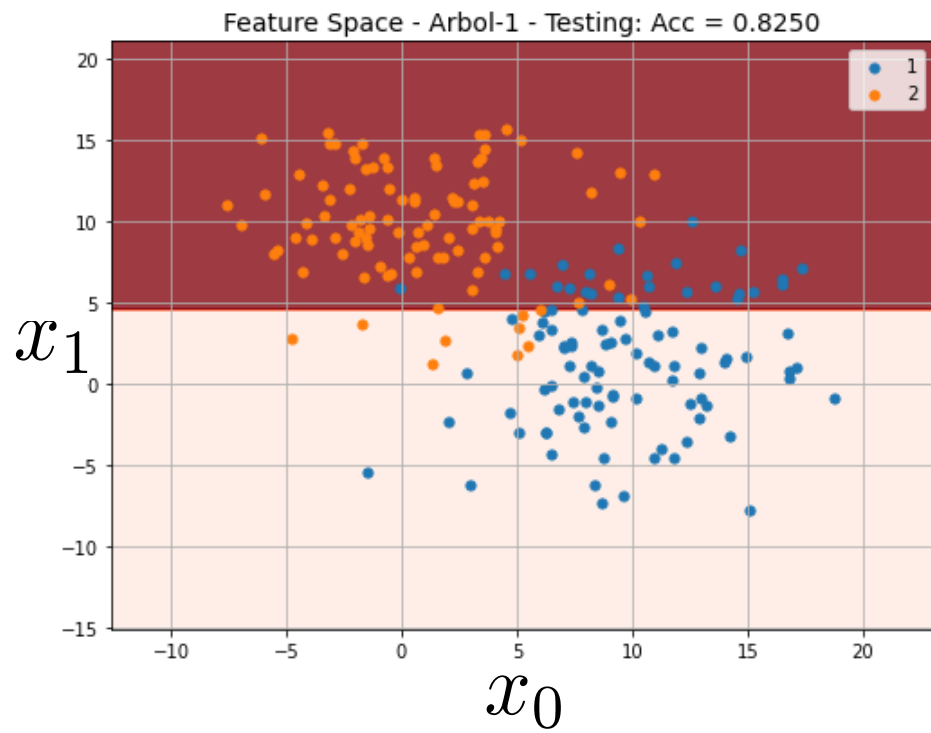
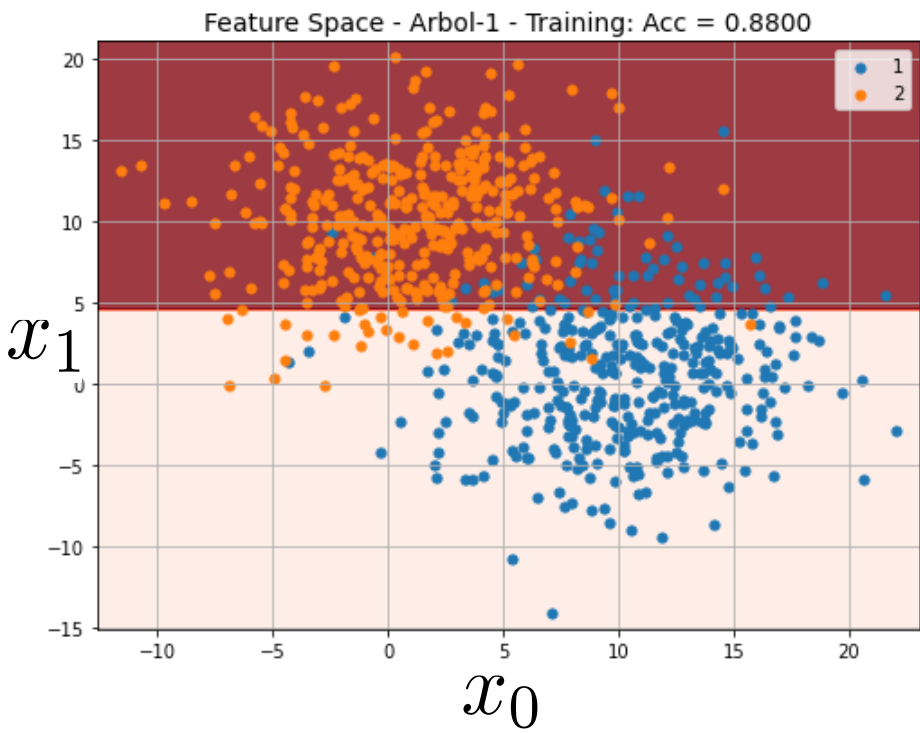


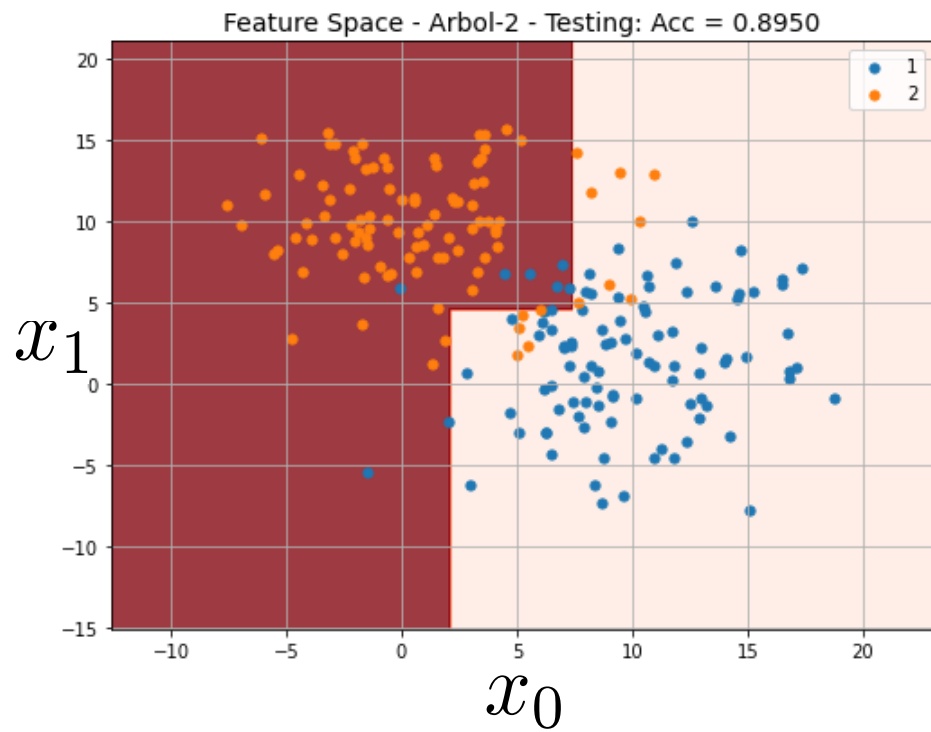
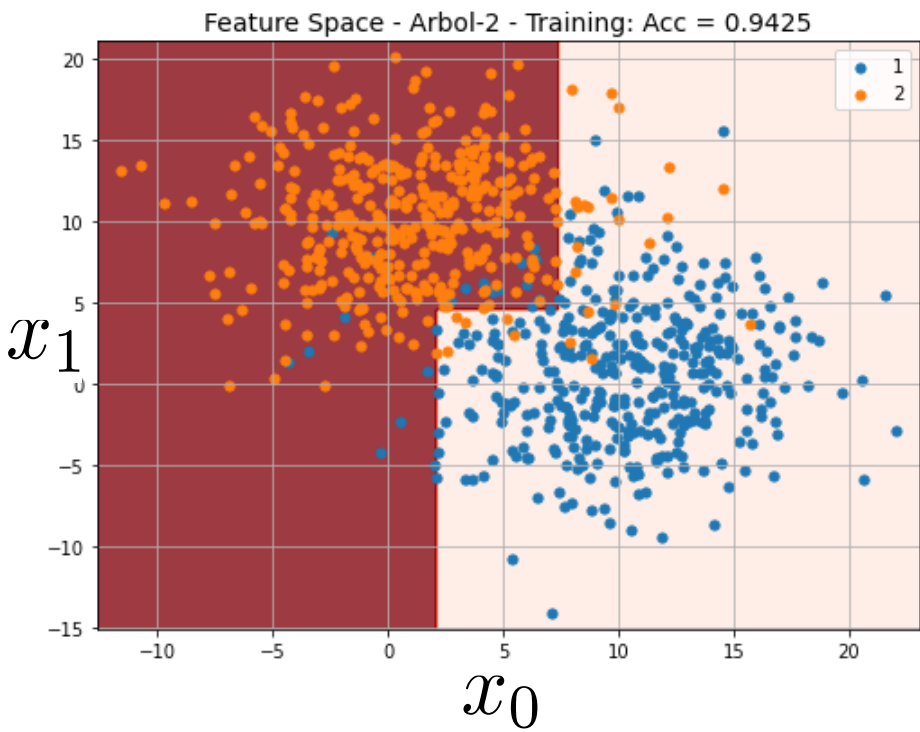
$X[1] \leq 4.508$
gini = 0.5
samples = 800
value = [400, 400]

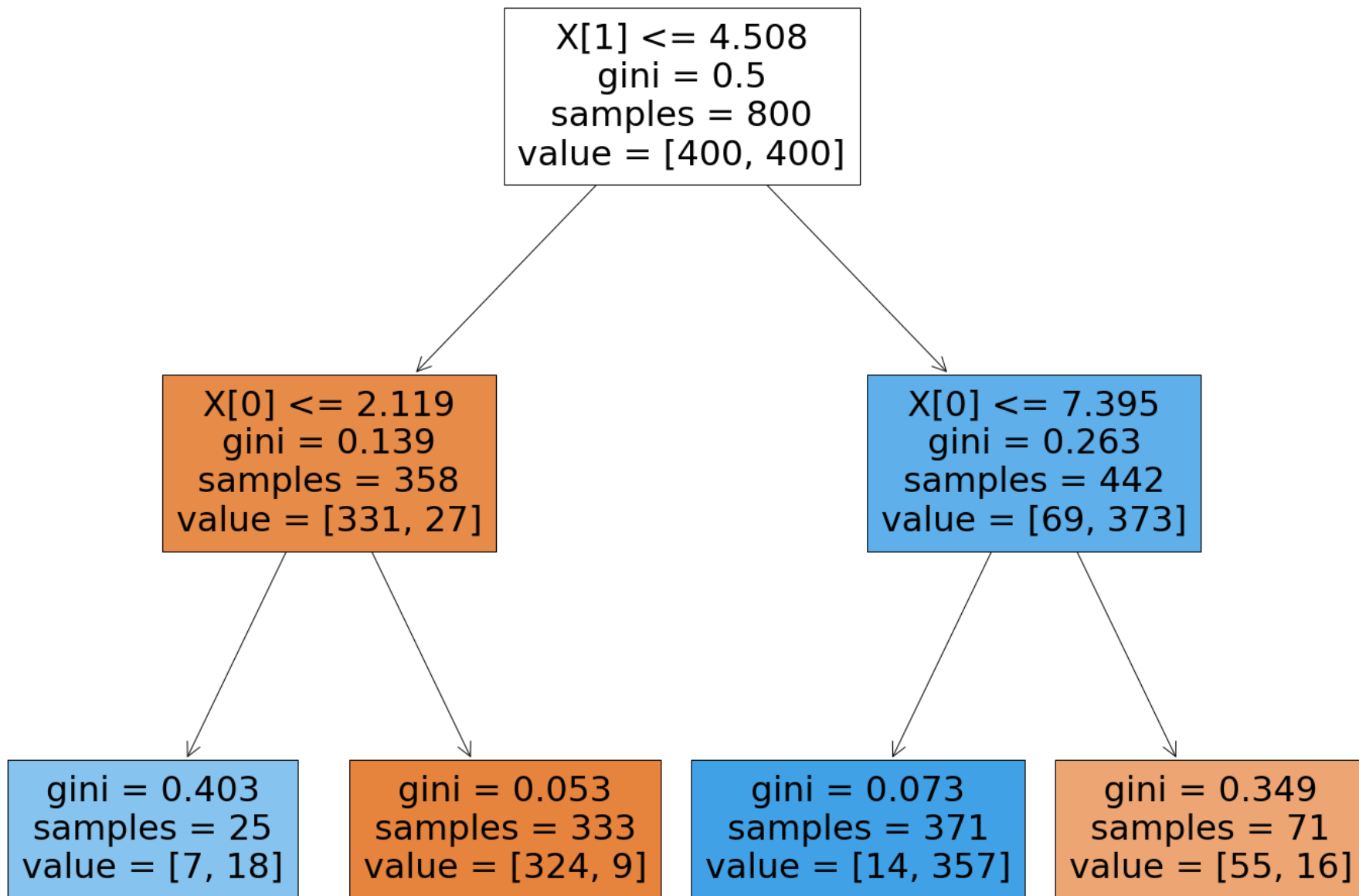


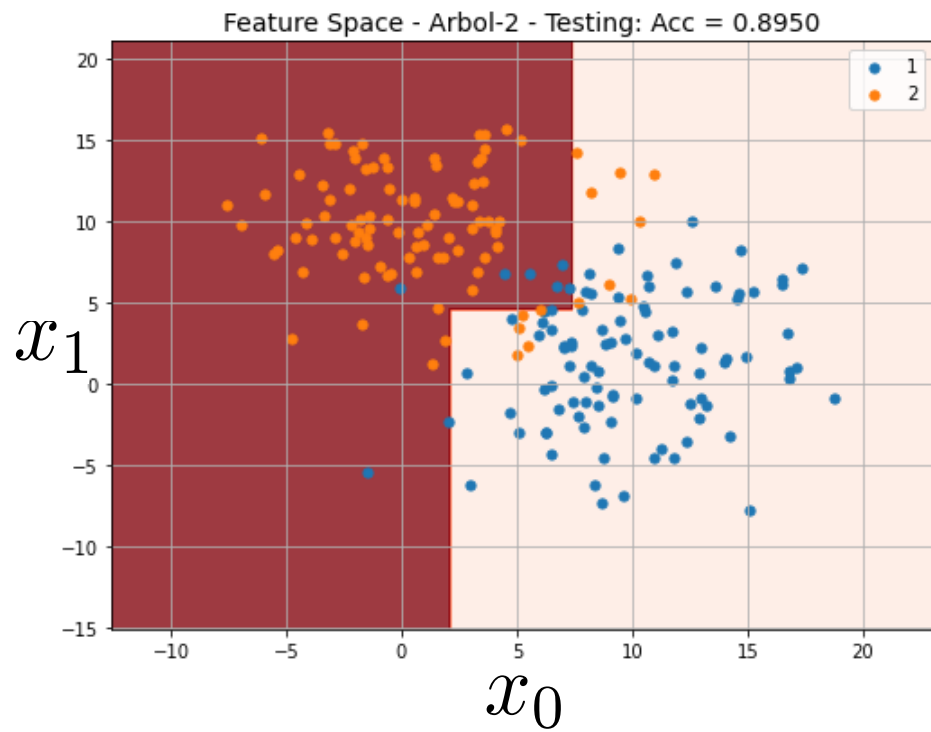
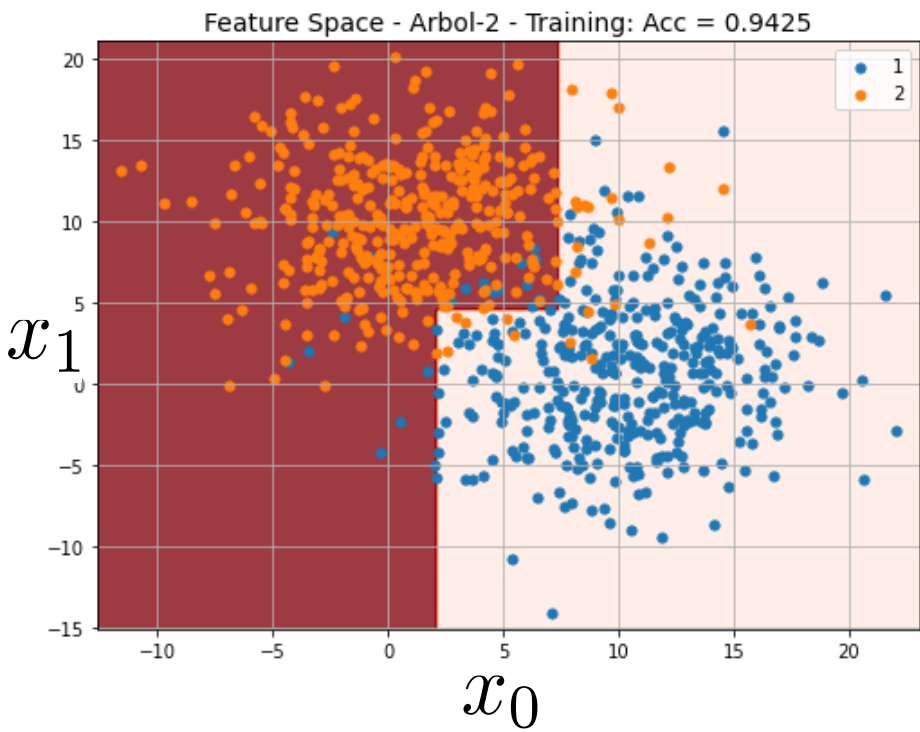
gini = 0.139
samples = 358
value = [331, 27]

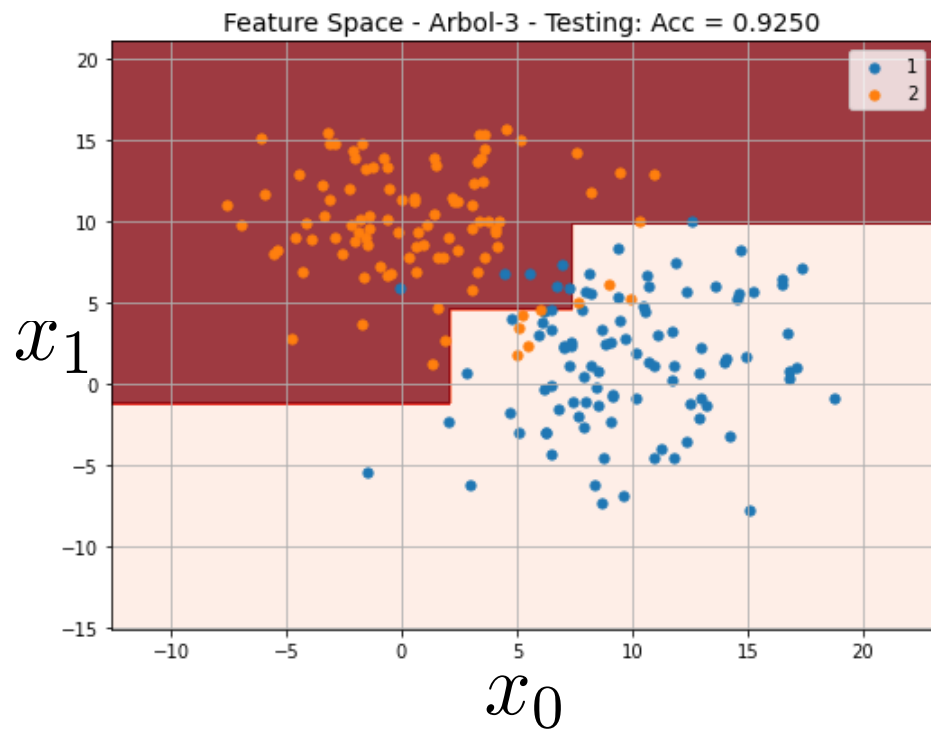
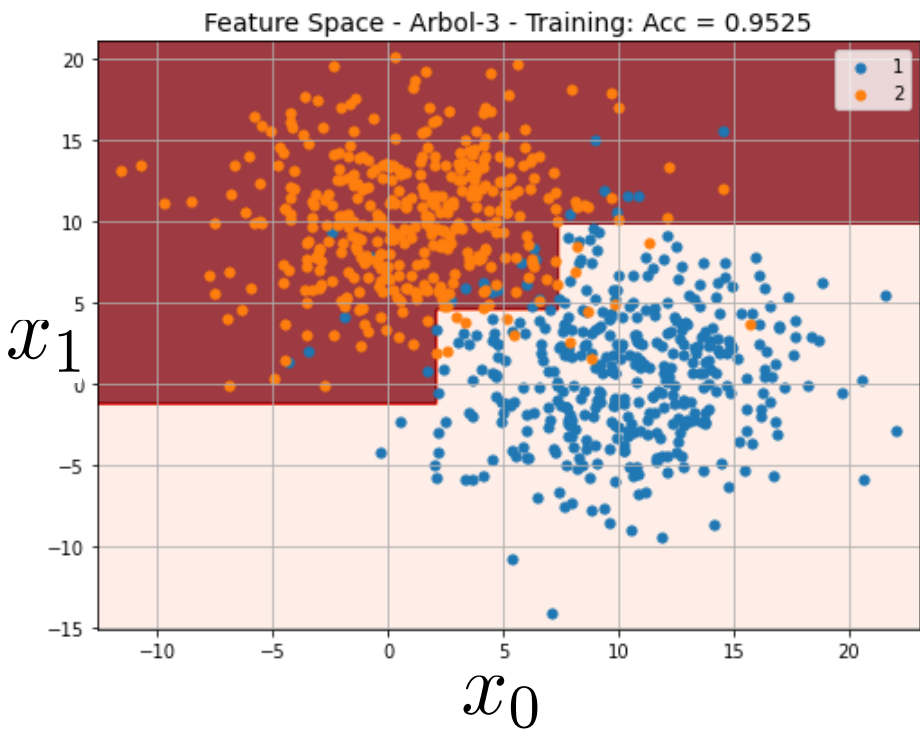
gini = 0.263
samples = 442
value = [69, 373]

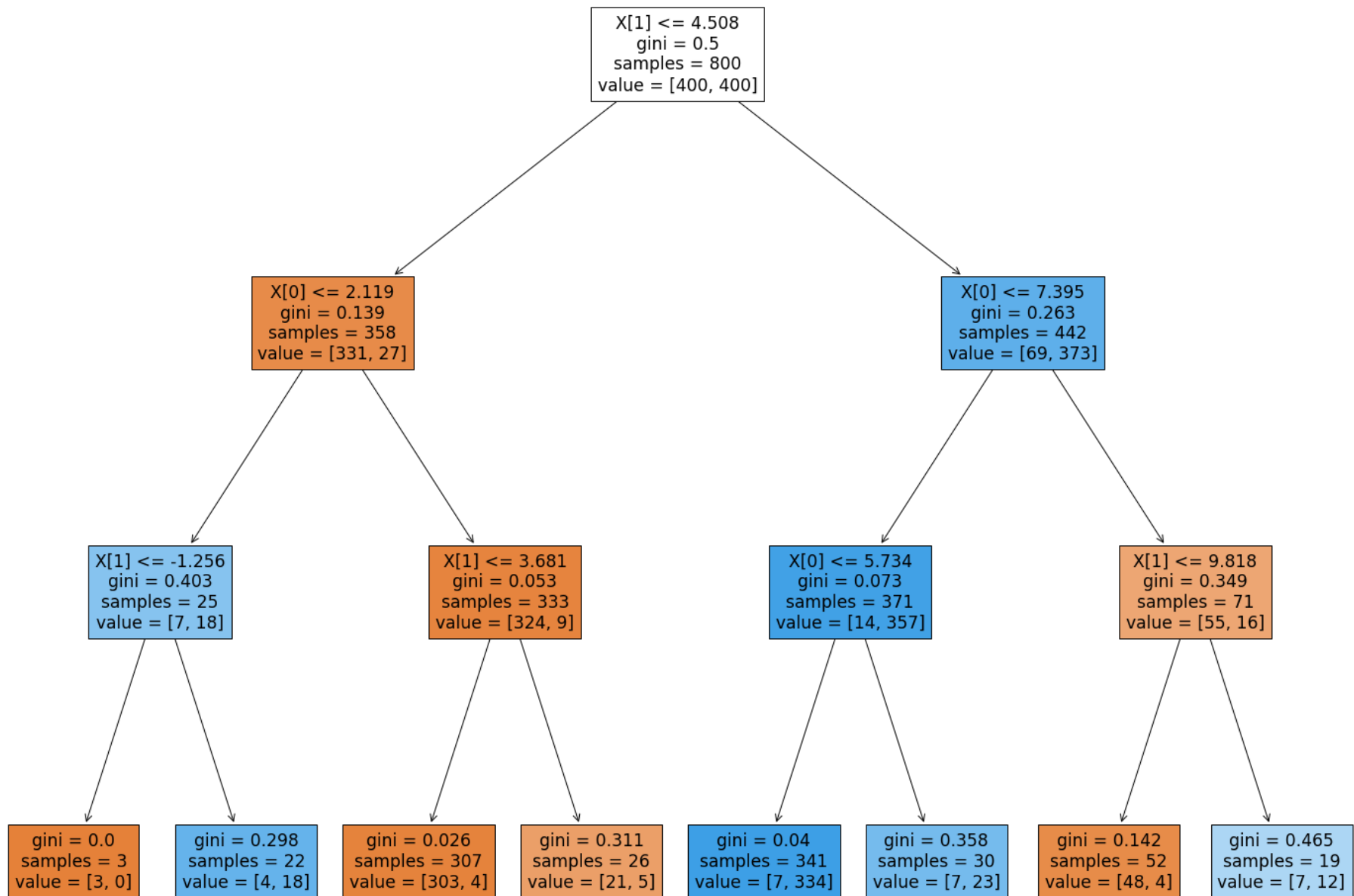


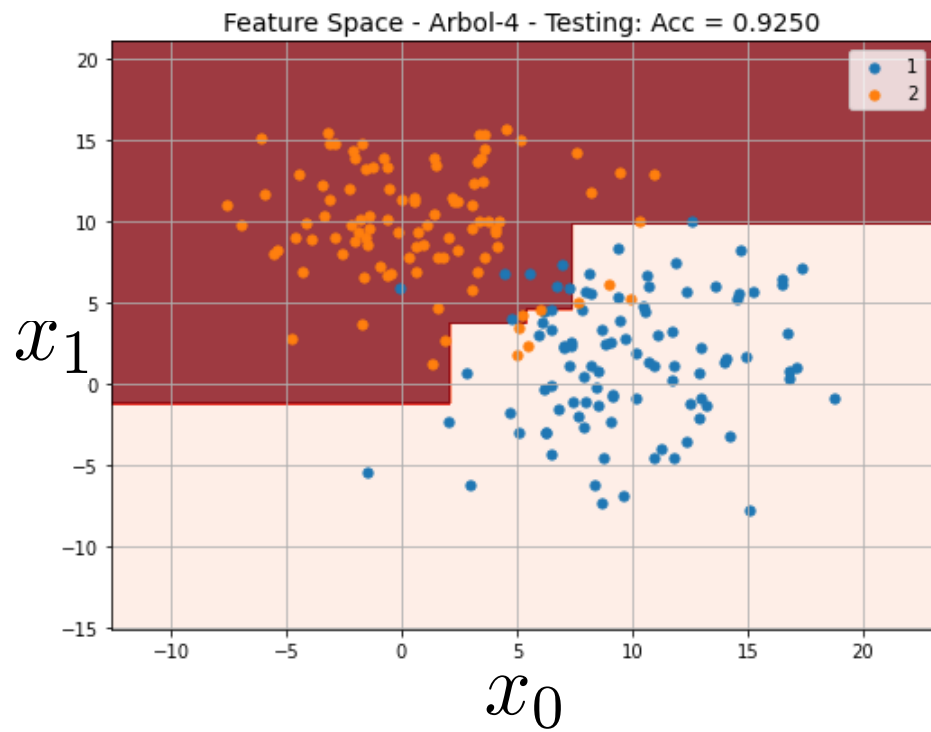
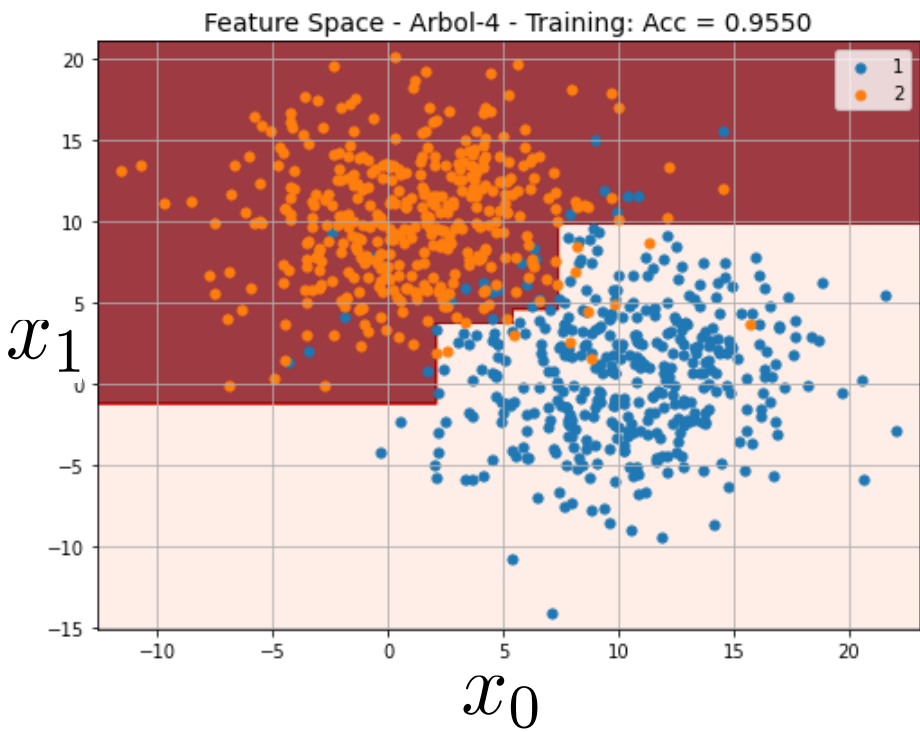


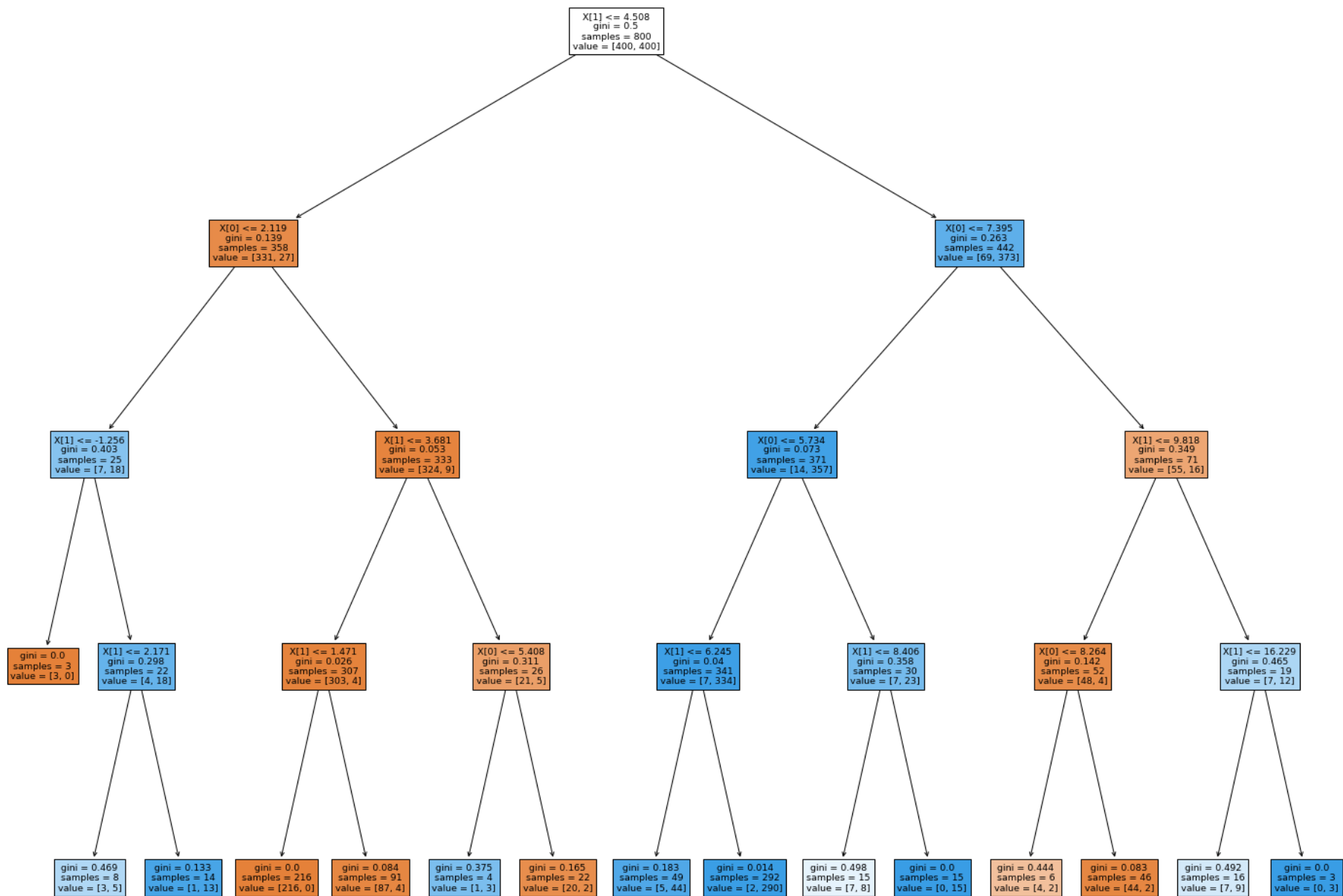


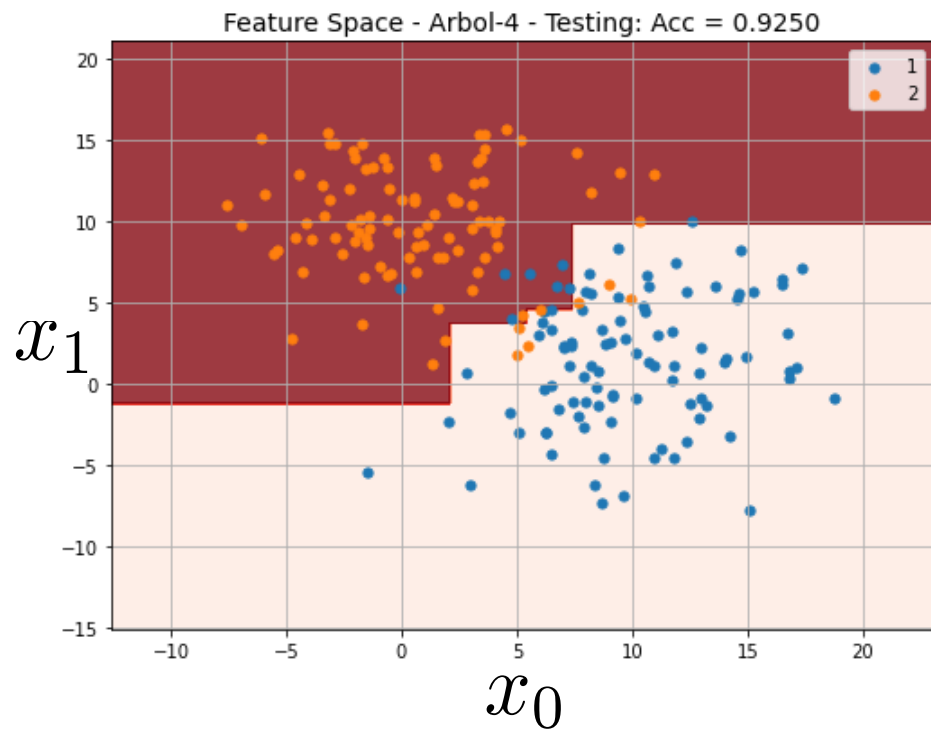
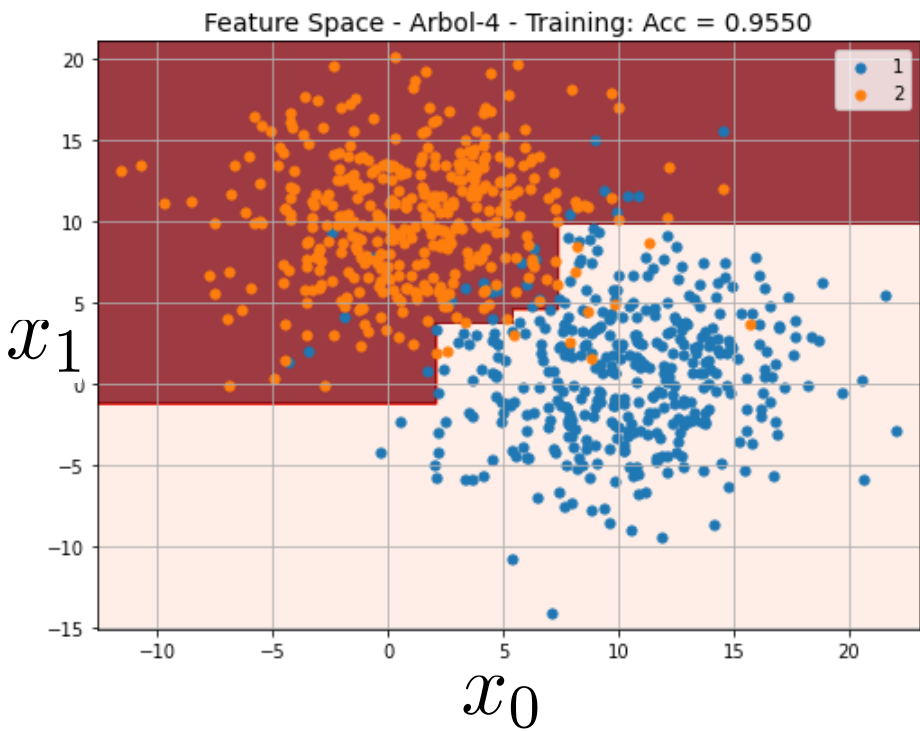


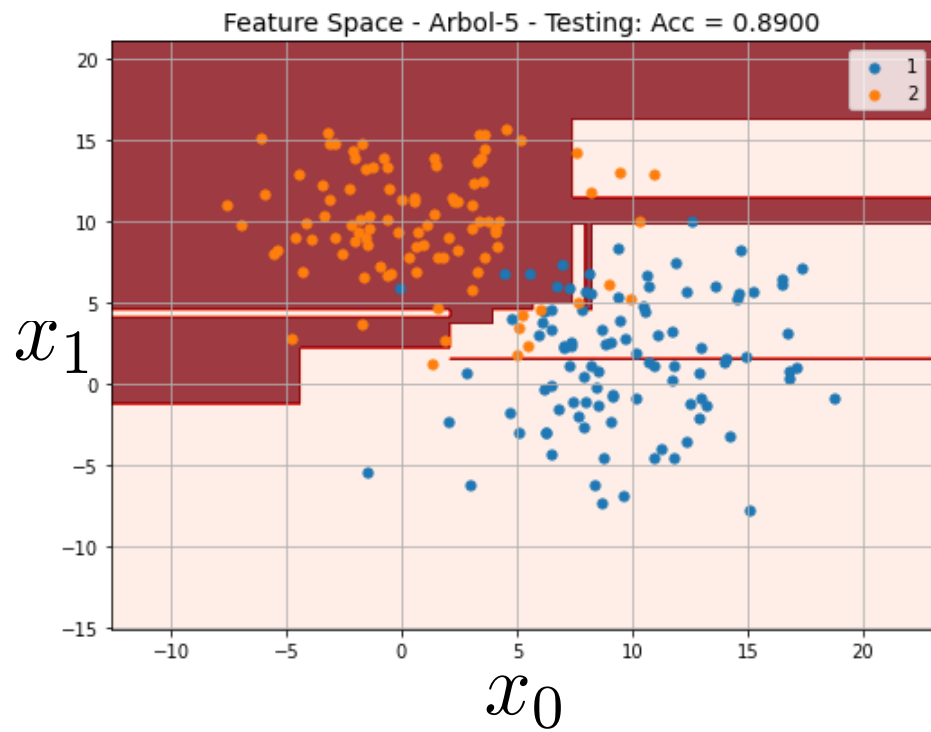
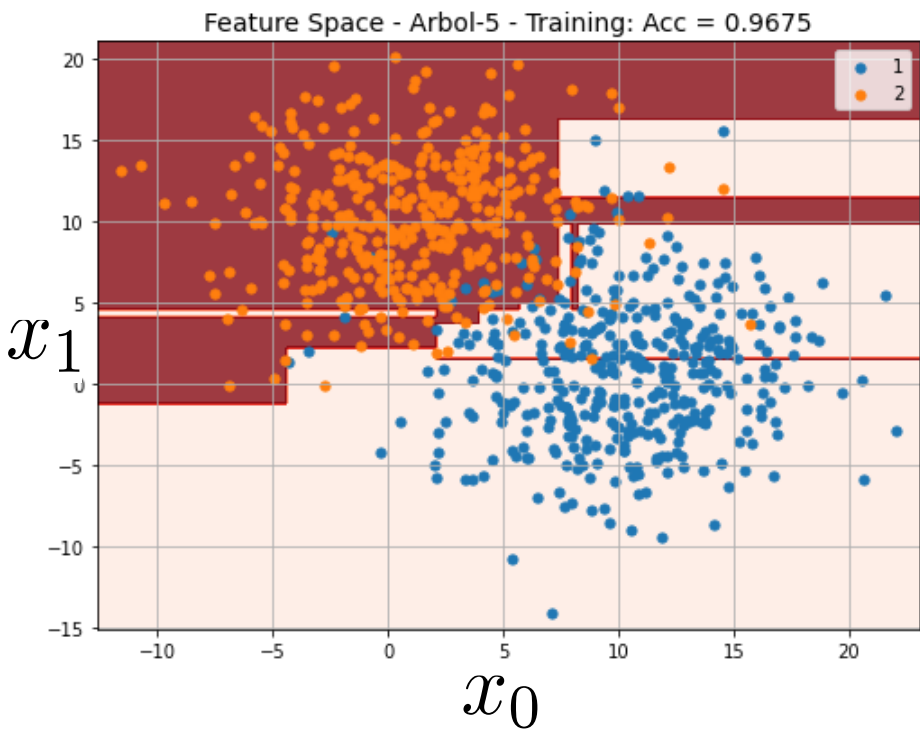


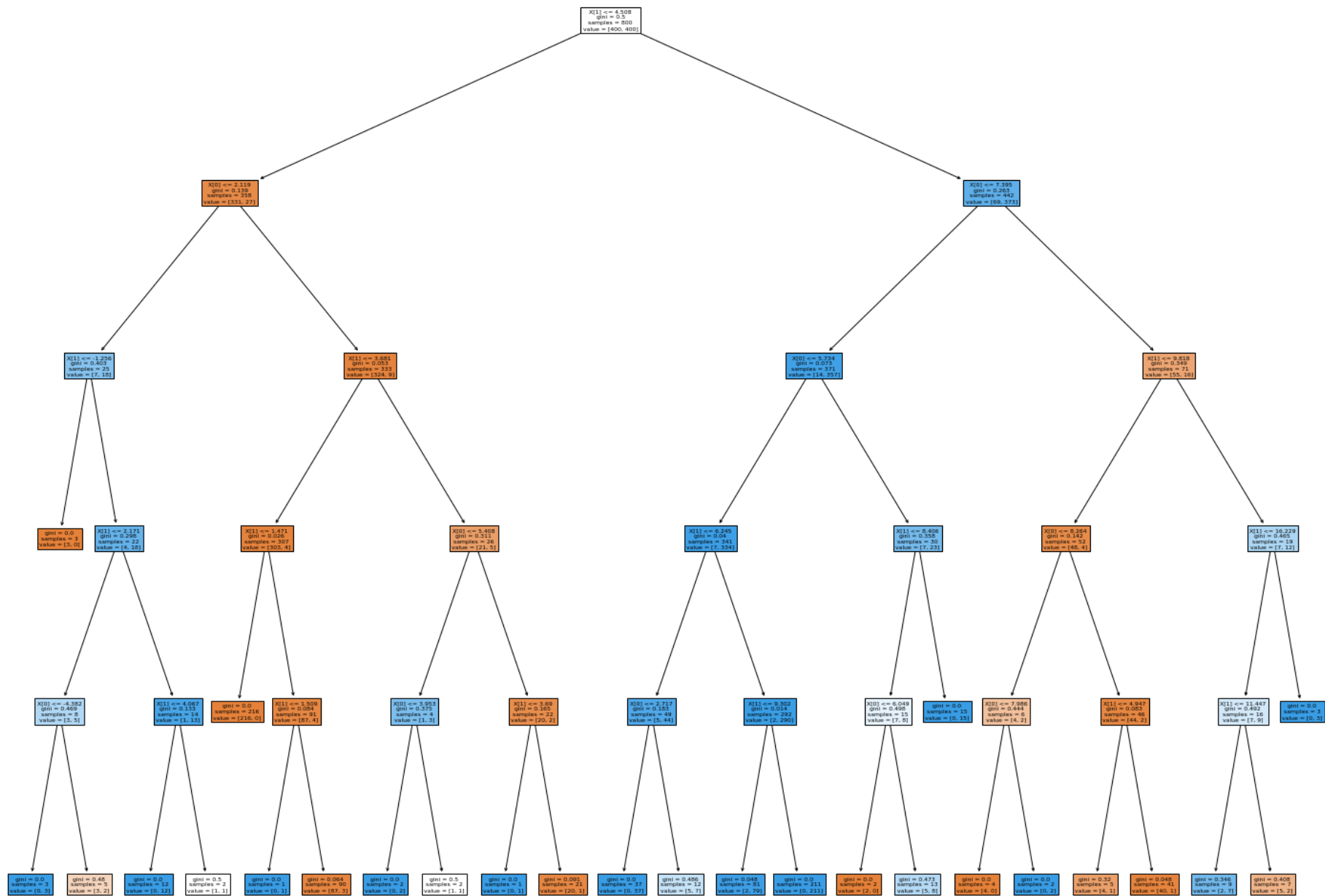












Random Forest

Random Forest

TRAINING:

for $i = 1$ to n

Escoger aleatoriamente un subconjunto de training

Entrenar un árbol de decisión A_i

Random Forest

TRAINING:

for $i = 1$ to n

Escoger aleatoriamente un subconjunto de training

Entrenar un árbol de decisión A_i

TESTING:

for $i = 1$ to n

Clasificar la muestra de testing usando A_i

Clasificar la muestra según la mayoría de los n votos

Random Forest

